# Video Concept Detection by Deep Nets with FLAIR

Cees Snoek, Koen van de Sande, Daniel Fontijne

Qualcomm Technologies
Netherlands B.V.

University of Amsterdam
The Netherlands

# Summary of our efforts

**Last year**

   Deep CNN for video concept detection and localization
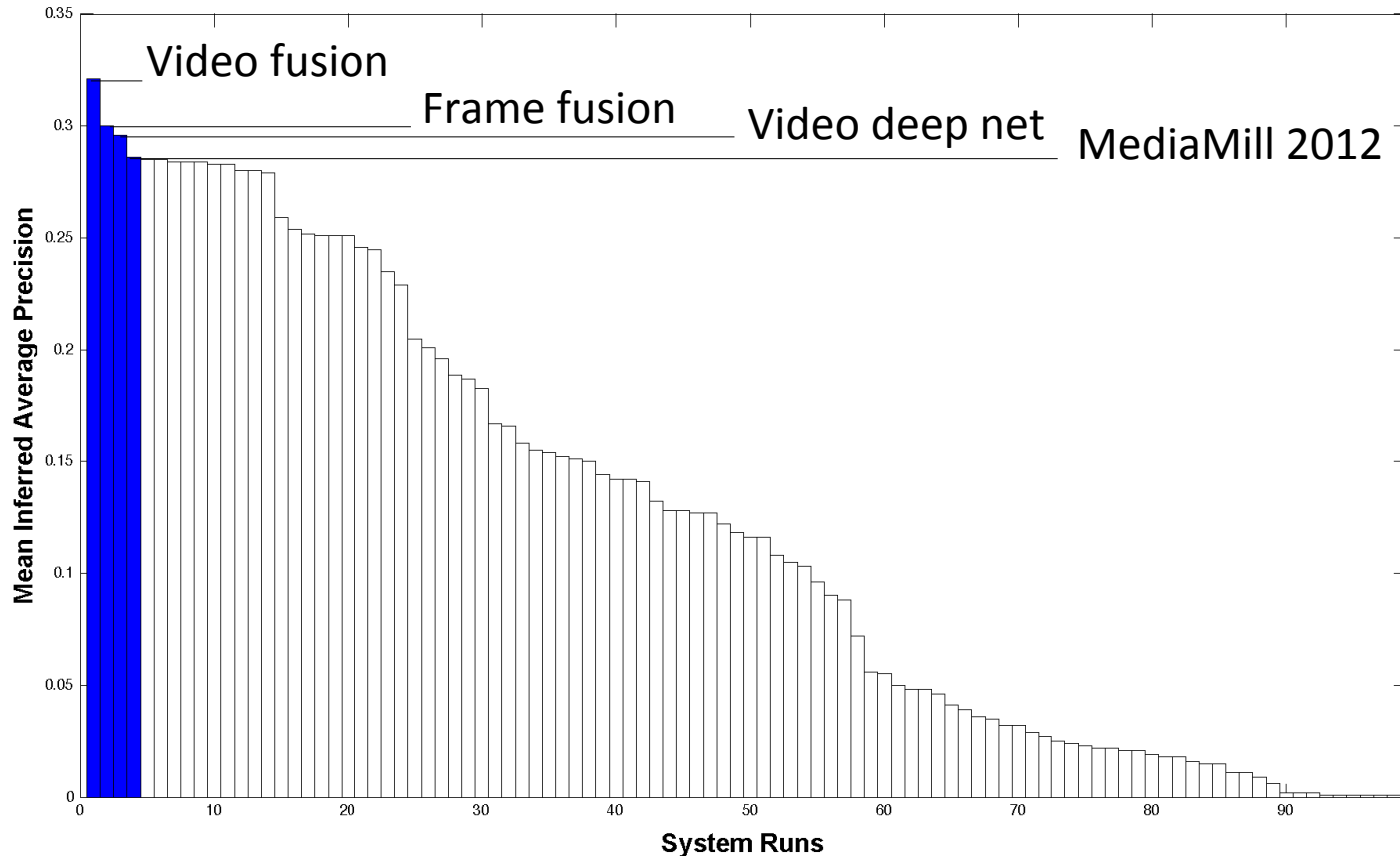
**This year**

   Tangential improvements for concept detection

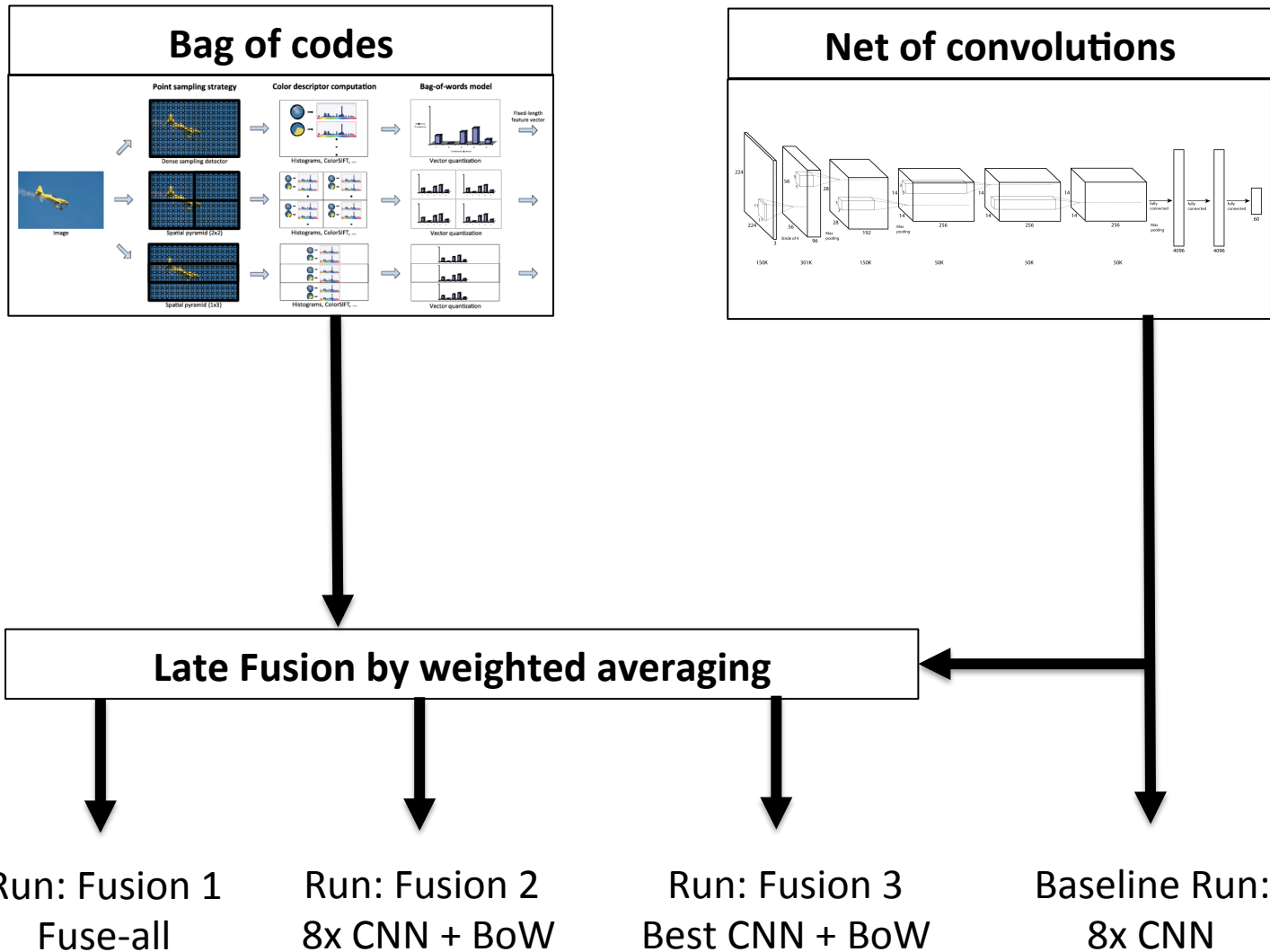   ***Our main innovation is in concept localization***
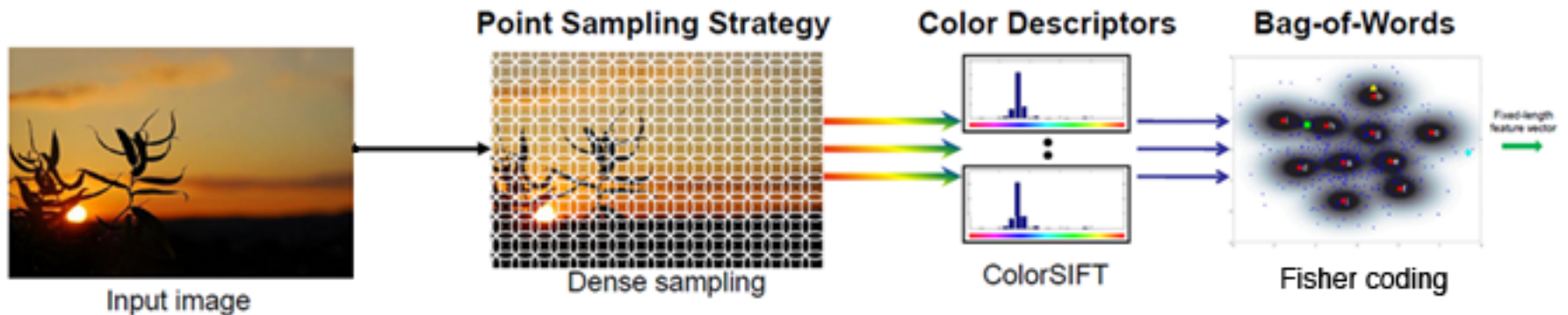
TASK I

**DETECTING CONCEPTS**

# Conclusion from TRECVID 2013



***Bag of words and deep net profit from each other***

# MediaMill TRECVID 2014 runs



**Bag of codes**

Point sampling strategy

Color descriptor computation

Bag-of-words model

Fixed-length feature vector

Dense sampling detector

Histograms, ColorSIFT, ...

Vector quantization

image

Spatial pyramid (2x2)

Histograms, ColorSIFT, ...

Vector quantization

Spatial pyramid (1x3)

Histograms, ColorSIFT, ...

Vector quantization

**Net of convolutions**

**Late Fusion by weighted averaging**

Run: Fusion 1
Fuse-all

Run: Fusion 2
8x CNN + BoW

Run: Fusion 3
Best CNN + BoW

Baseline Run:
8x CNN

# MediaMill: Color difference coding

- Densely sampled points
- SIFT, C-SIFT and T-SIFT descriptors
- PCA reduction to 80D
- Fisher vector coding with codebook size 256
- Spatial pyramid 1x1+1x3
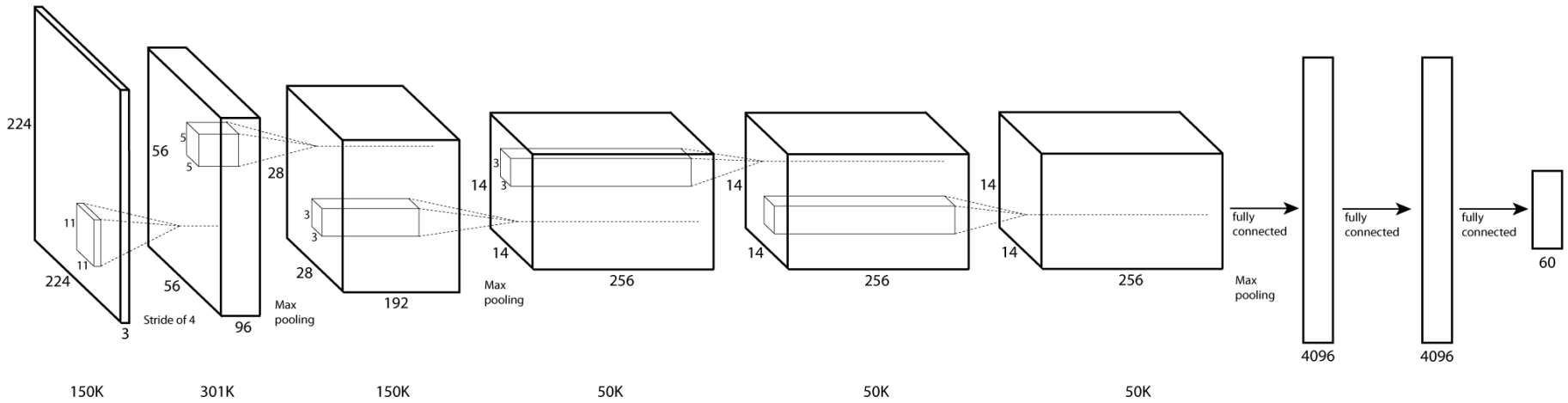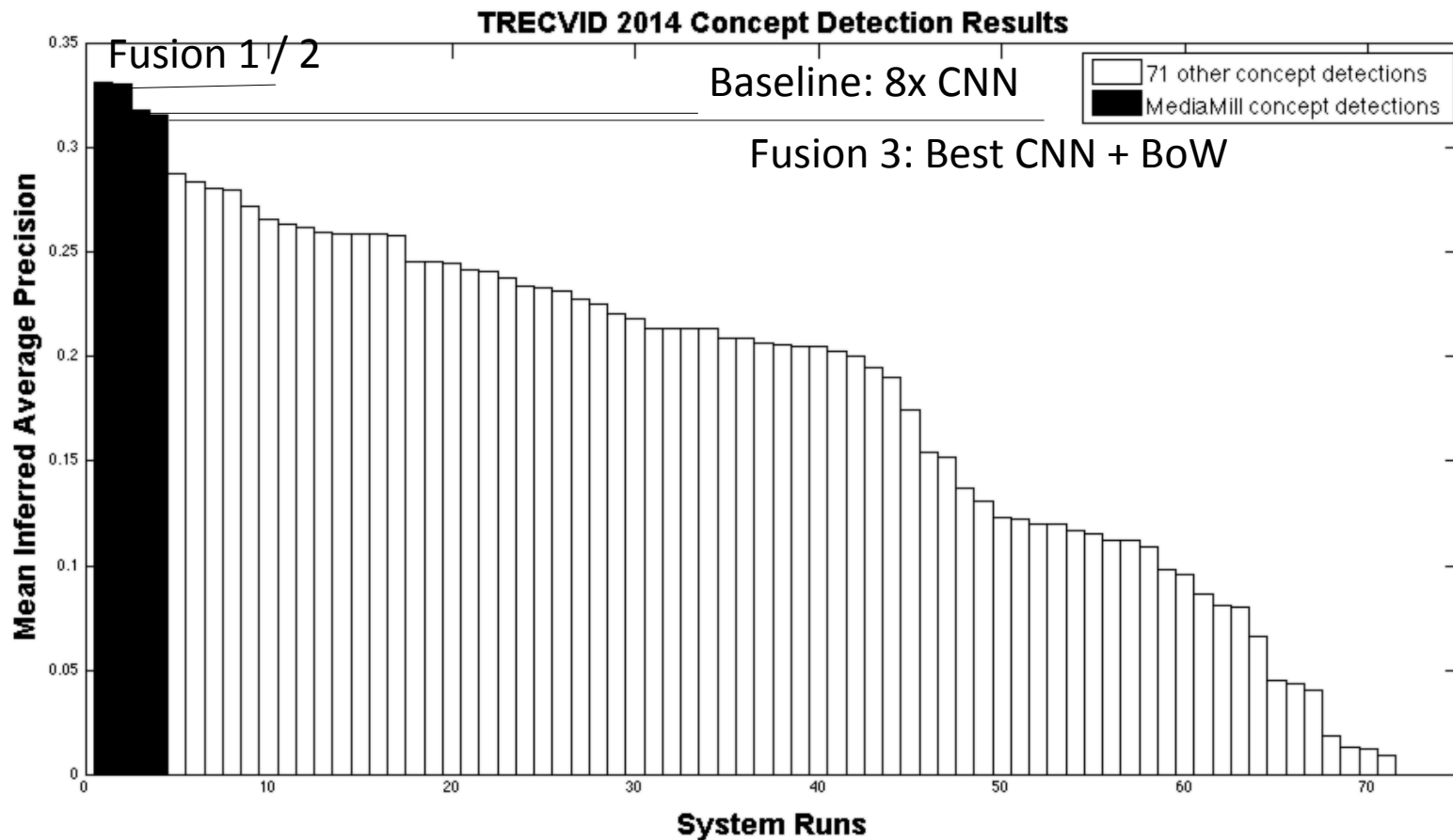- Spatial coordinate coding
- Linear classifier



**Point Sampling Strategy** — Dense sampling

**Color Descriptors** — ColorSIFT

**Bag-of-Words** — Fisher coding

Input image

Fixed-length feature vector

# MediaMill: Video deep learning

Convolutional neural network with 8 layers with weights

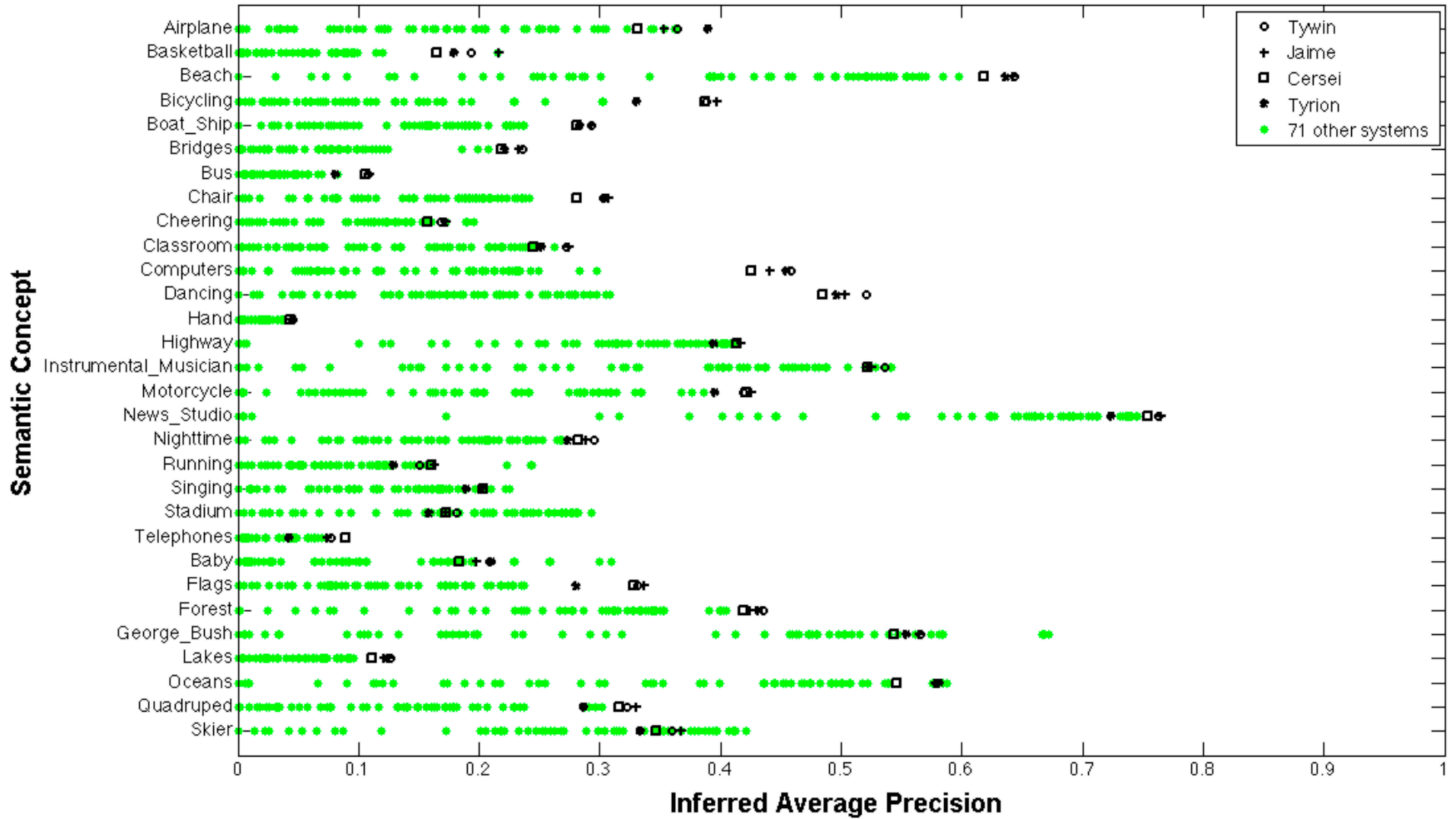Trained using error back propagation
- ImageNet for pre-training

# Results



**TRECVID 2014 Concept Detection Results**

Fusion 1 / 2

Baseline: 8x CNN

Fusion 3: Best CNN + BoW

71 other concept detections
MediaMill concept detections

Mean Inferred Average Precision

System Runs

*Bag of words and deep net profit from each other, better results with more nets*

# Results per concept



TRECVID 2014 Semantic Indexing Task Benchmark Comparison

# LOCALIZING CONCEPTS

# Goal: meaningful localization



Finding **where**, **when**, **what** is happening

Challenges: huge search space, non-rigid deformation

# Related work

## Sliding Window

Image                     Video



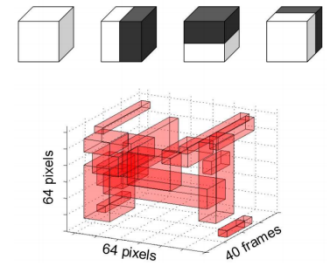[Rowley, 1996]            [Rodriguez, 2008]

## Boosting Cascade

Image                     Video



[Viola & Jones, 2001]     [Ke, 2005]

## Branch and Bound

Image                     Video
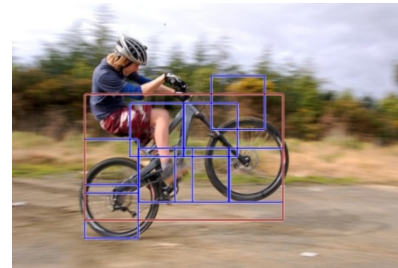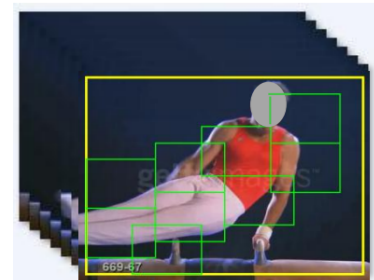


[Lampert, 2009]           [Yuan, 2011]

## Deformable Parts

Image                     Video
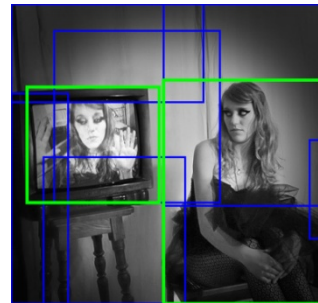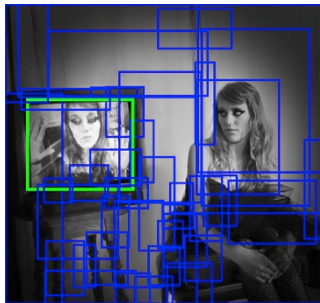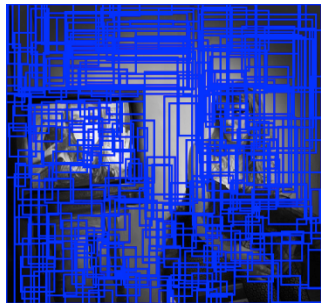


[Felzenswalb, 2008]       [Tian, 2013]

# Inspiration: Selective Search

[Uijlings, 2013]

Iterations of selective search →



Hierarchical grouping of super-pixels
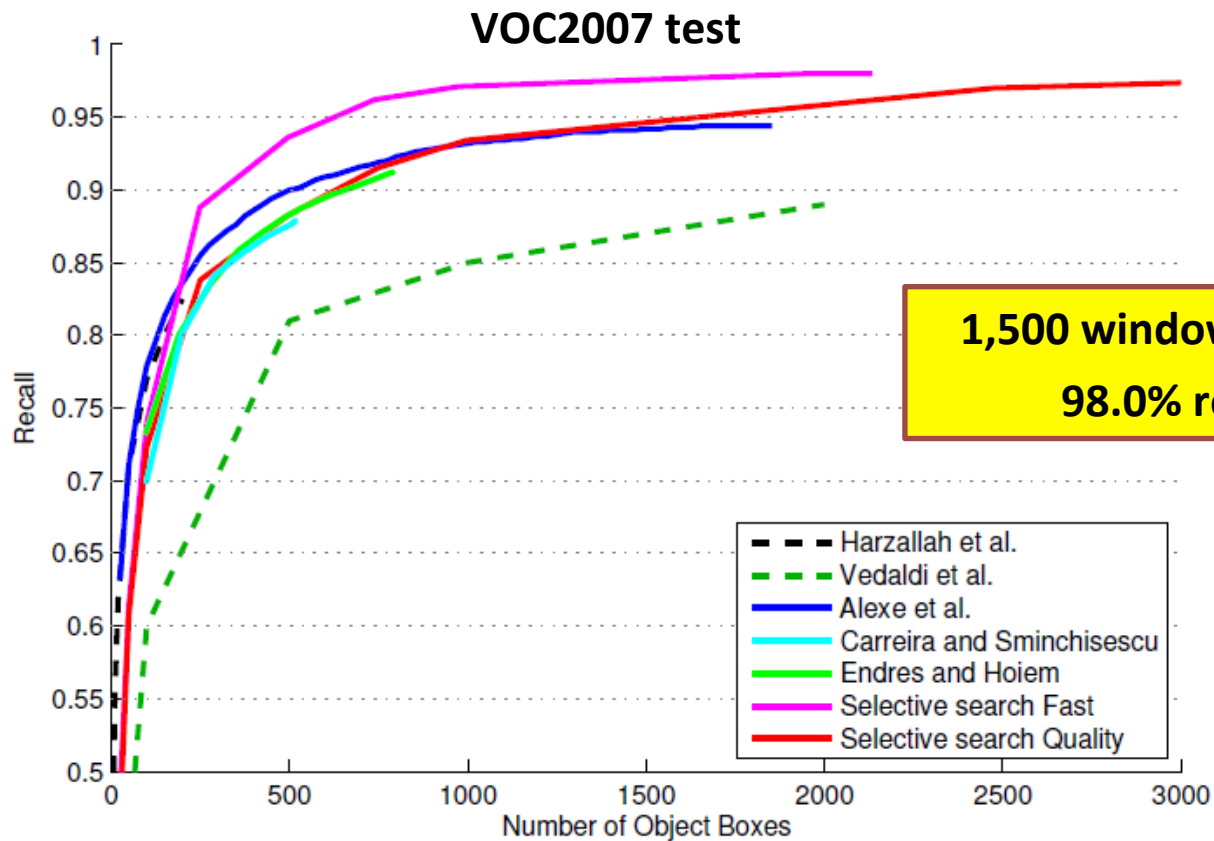
Object proposals

High recall with modestly sized  object hypotheses set

Feasible to train an expensive classifier

# Selective Search

Multiple complementary invariant color spaces
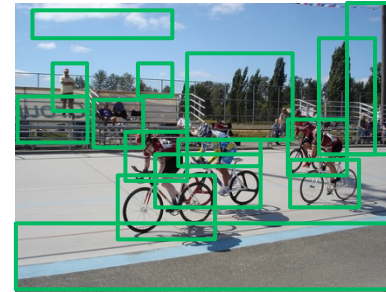
Location hypotheses are class-independent

**VOC2007 test**



**1,500 windows/image**

**98.0% recall**

**Software available for download at http://koen.me/research/selectivesearch/**
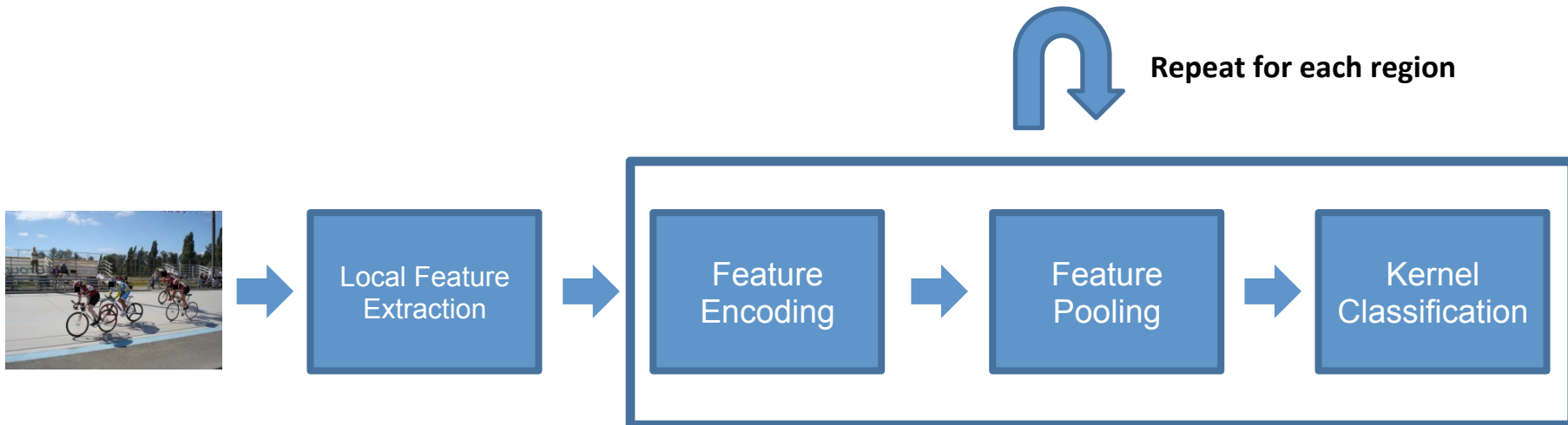
# Local object classification

Requires **repetitive** computations on **overlapping** regions

**Spatial Pyramids** [*Lazebnik, CVPR06*]
(#regions: 10-100)

**Object Detection** [*Sande, ICCV11*]
(#regions: 1,000-10,000)

**Repeat for each region**

| Local Feature Extraction | → | Feature Encoding | → | Feature Pooling | → | Kernel Classification |

# Features

Use SIFT and ColorSIFT descriptors

Bag-of-words, VLAD, Fisher vector encoding

Encoding 2000 boxes per image is expensive

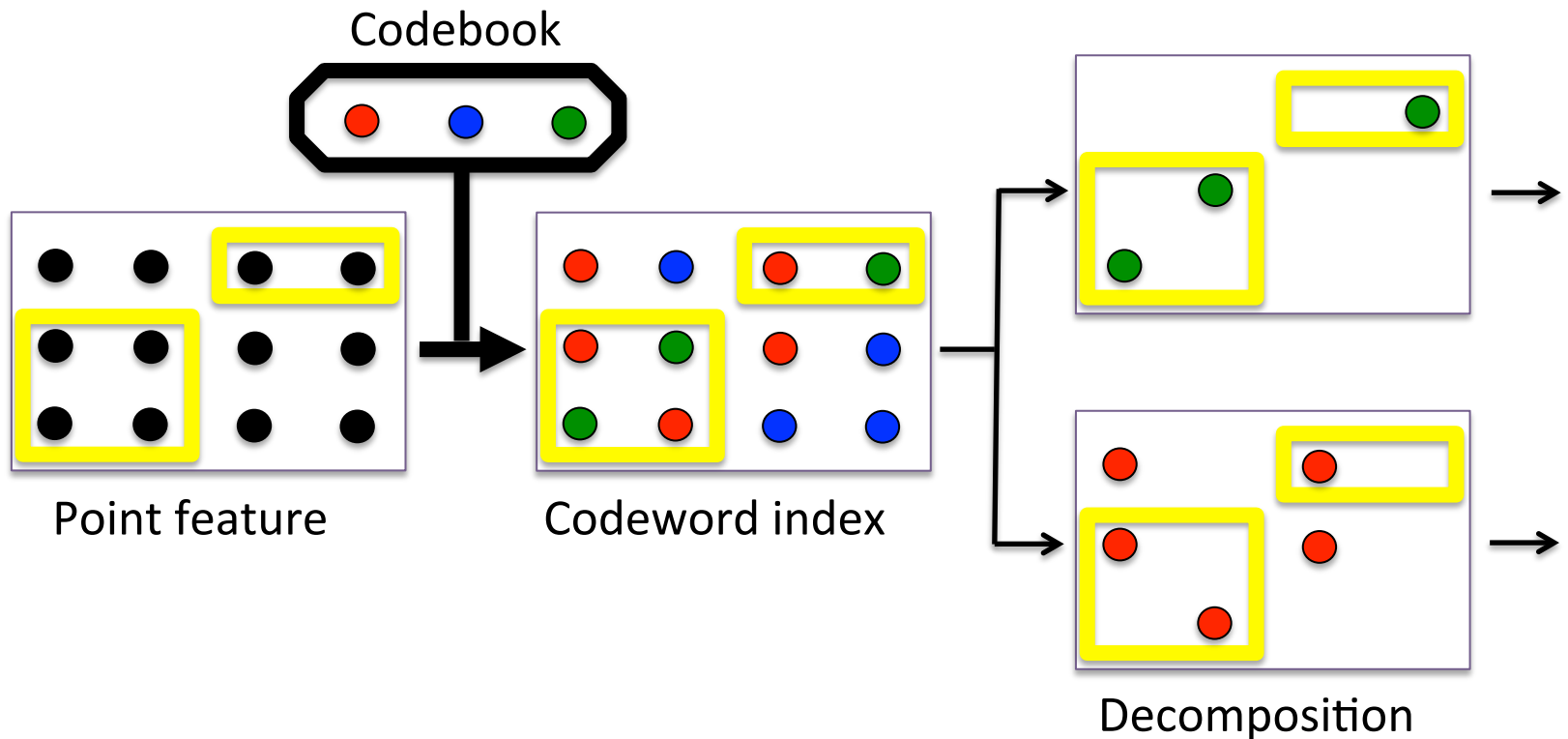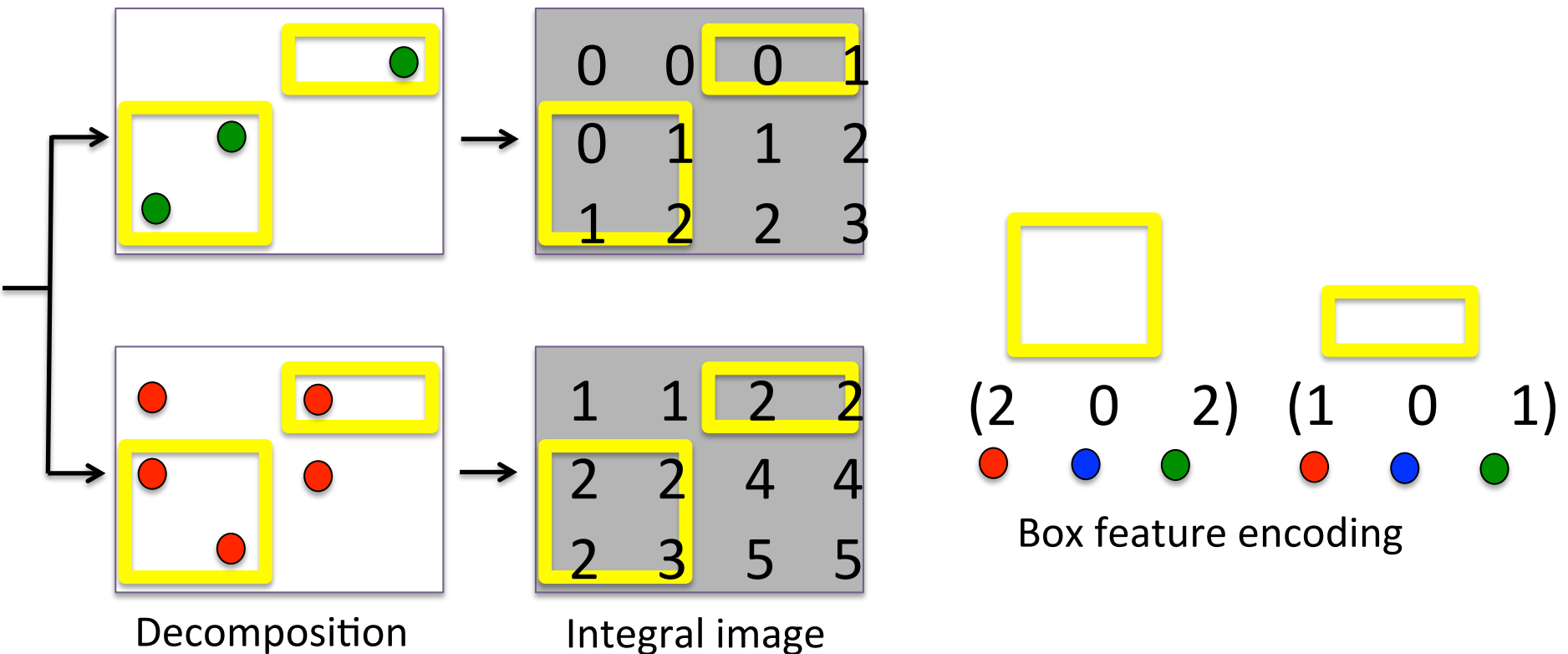    Bag-of-words:    10s

    VLAD:    30s

    Fisher:    120s

# Key idea

Decompose assignment over codebook elements

Codebook

Point feature

Codeword index

Decomposition

# Area-independent decomposition

Fast box evaluation with integral images



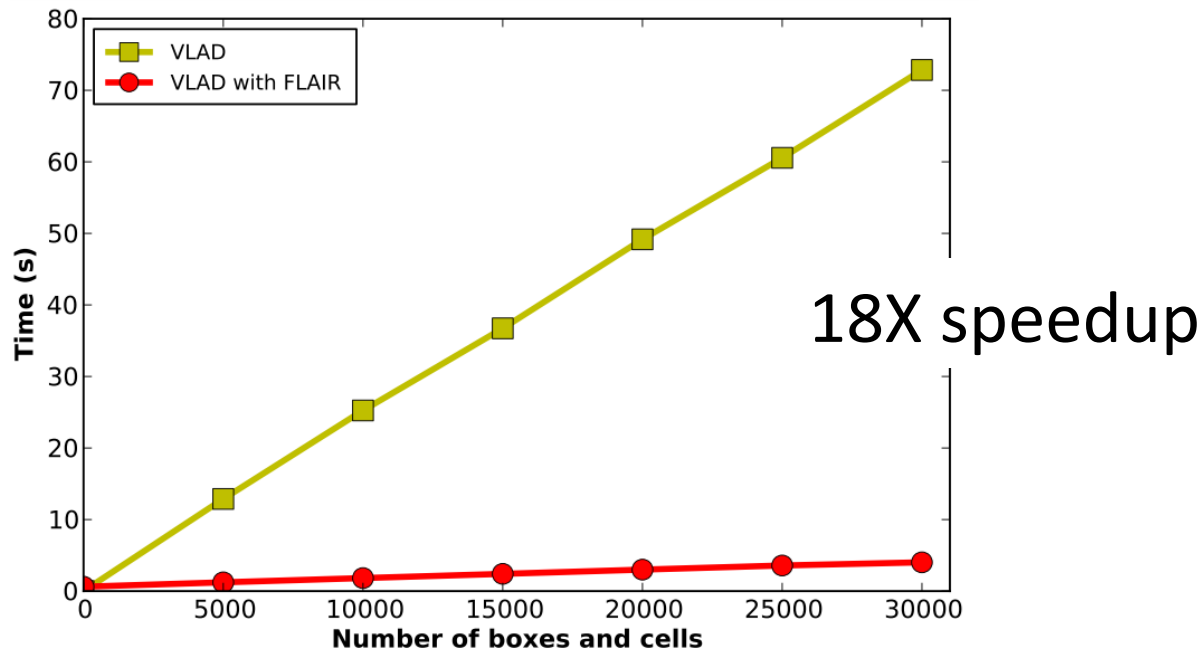Decomposition    Integral image

Box feature encoding

# VLAD with FLAIR

Decomposition as multi-dimensional integral image

Sparsity drops memory from 14GB to 1GB/image
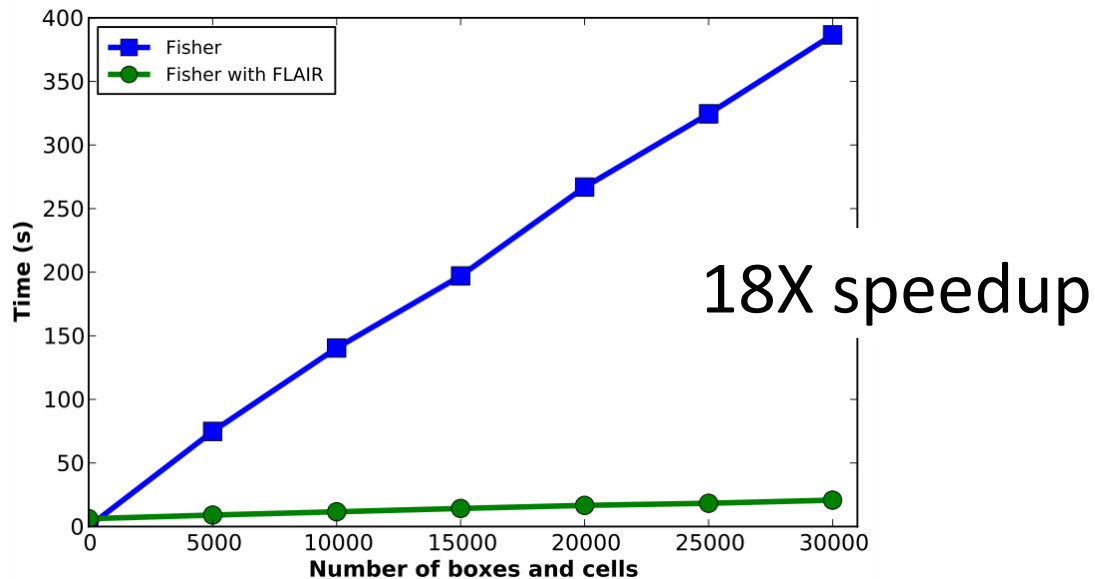
Supports power norm, L2 norm and spatial pyramid



18X speedup

# Fisher with FLAIR

Decomposition as four multi-dimensional integral images [See paper]

Supports power norm, L2 norm, spatial pyramids

No need for approximations

Scalable to modern datasets

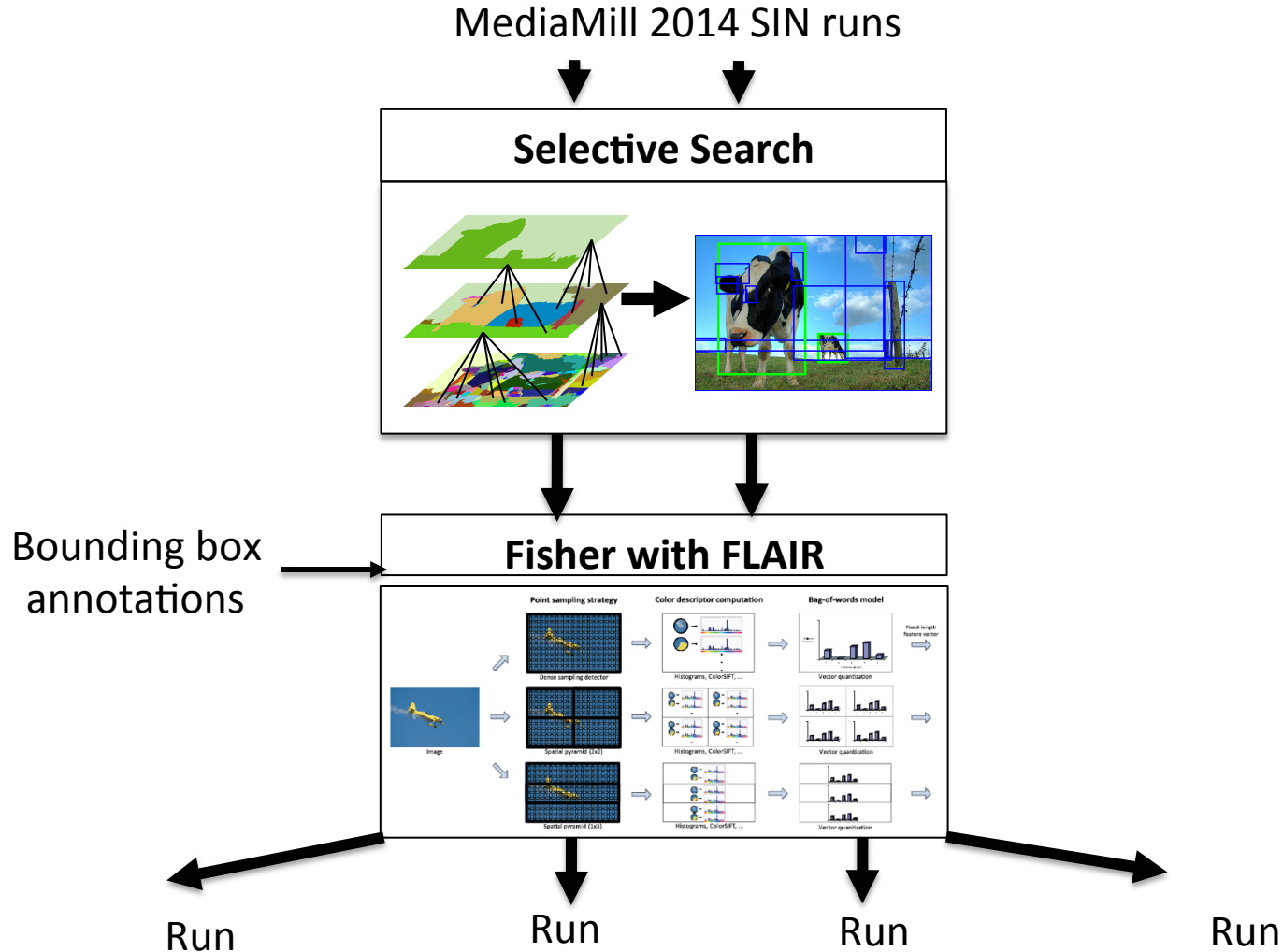18X speedup

# Overall detection speedup and accuracy

| | Time (s) per image | | | |
| --- | --- | --- | --- | --- |
| | Standard | with FLAIR | Speedup | mAP |
| BoW | 47.9 | - | - | 32.3 |
| VLAD | 34.3 | 7.8 | 4.4x | 28.2 |
| Fisher | 120.0 | 32.5 | 3.7x | 33.3 |

*Fisher with FLAIR is better and faster than BoW*

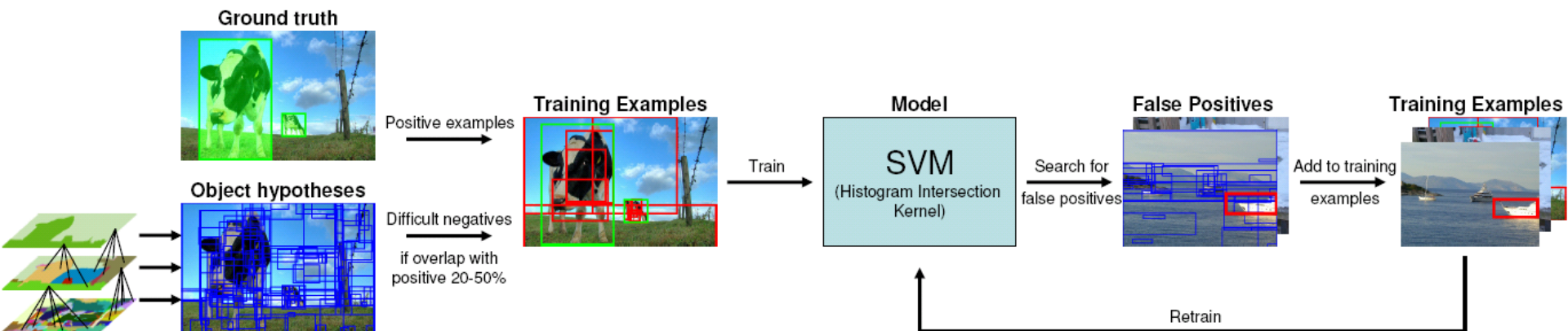# MediaMill TRECVID 2014 runs

# Implementation details

PCA-reduced ColorSIFT descriptors to 80D

Fisher with FLAIR encoding

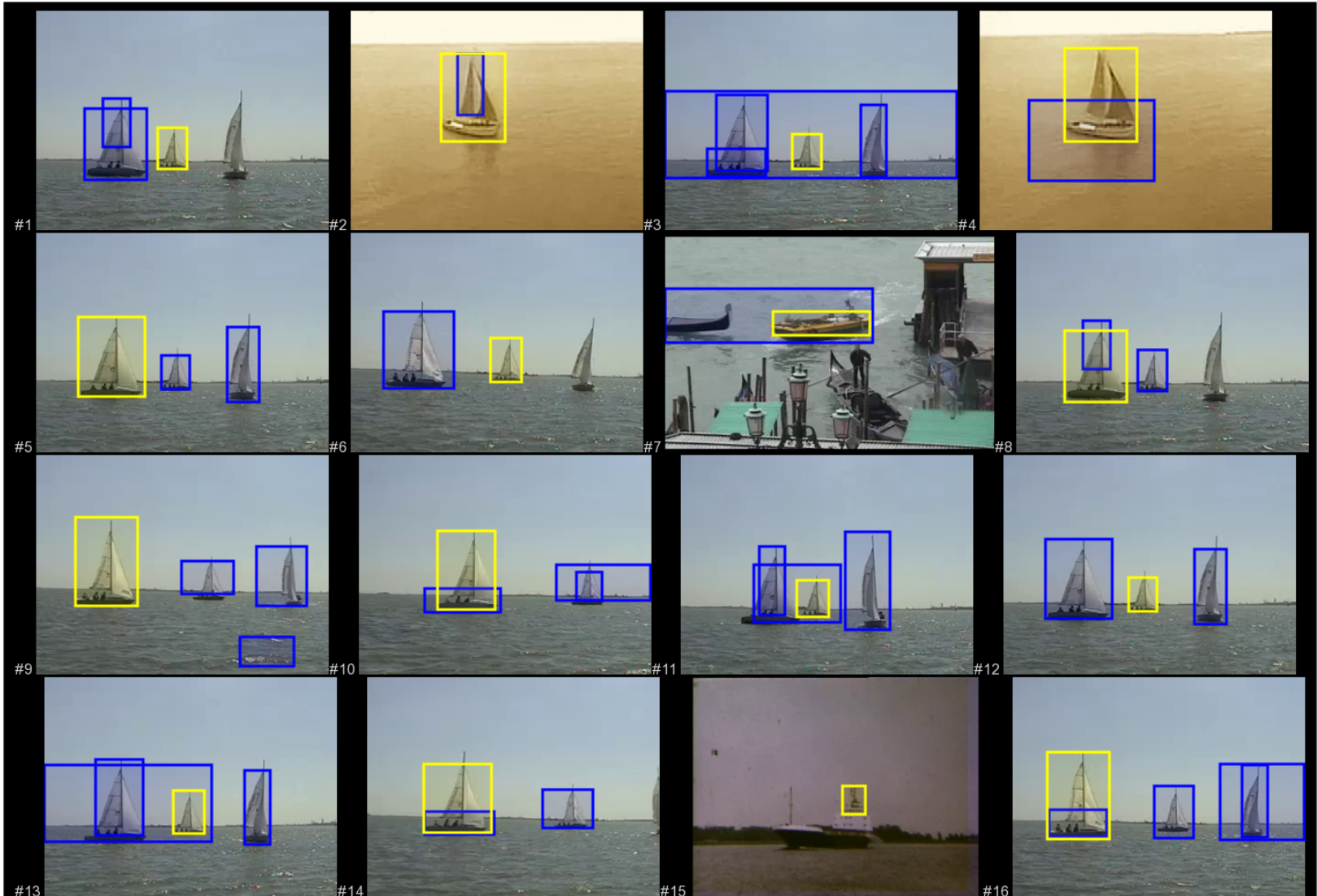Spatial pyramid

Linear SVM

Hard negative mining
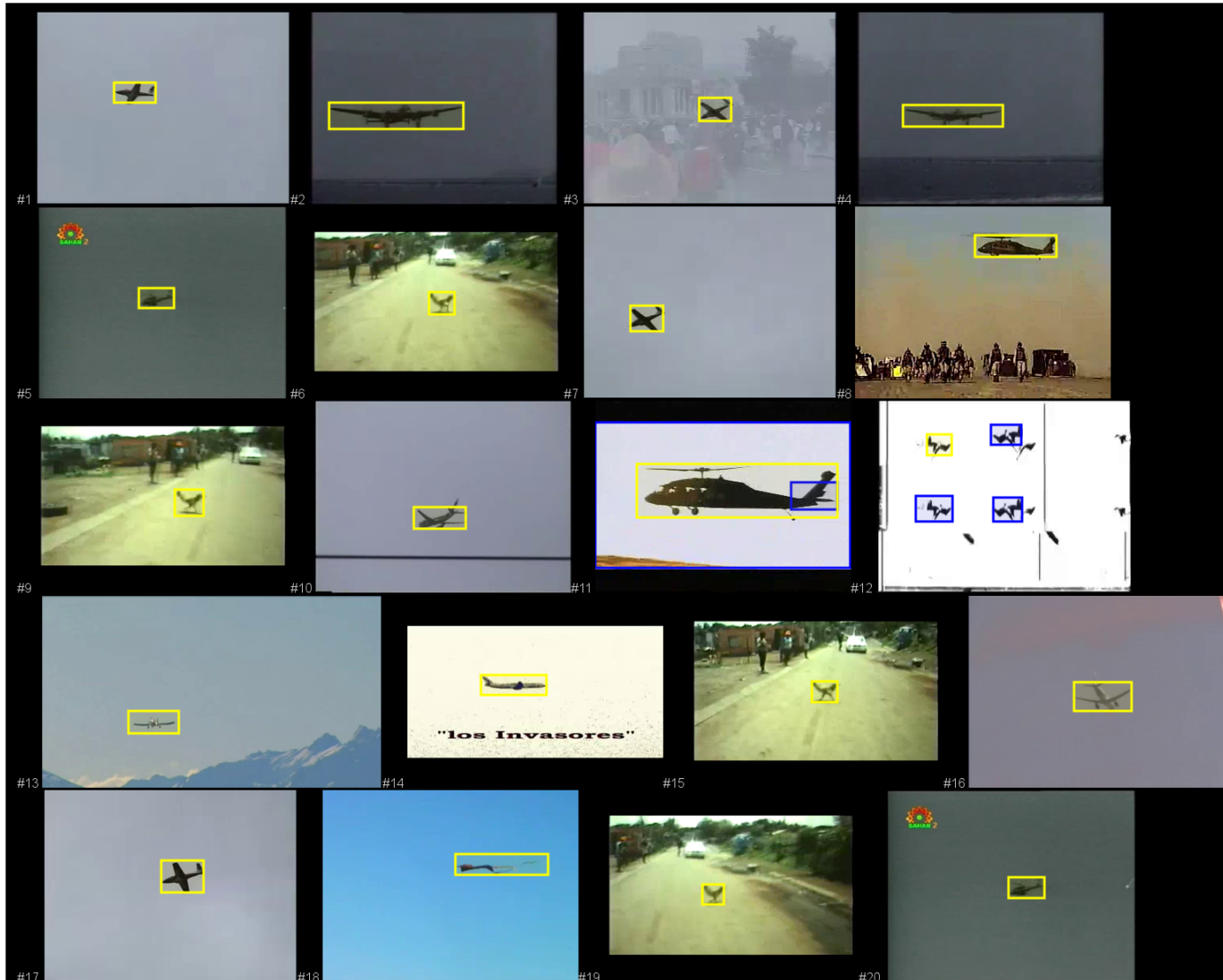
# Boat
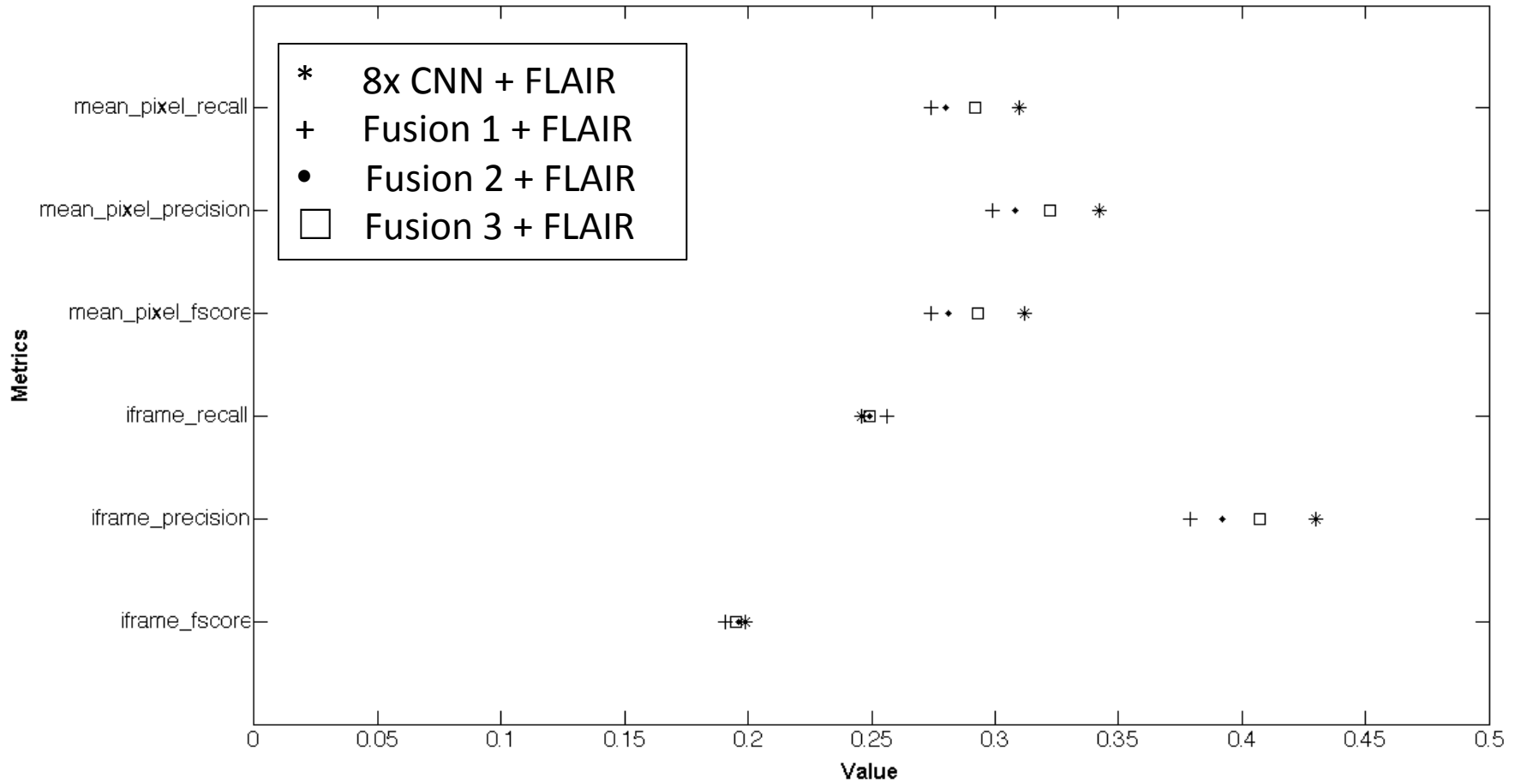
Best box

Other boxes

# Airplane

# Results



***FLAIR after deep nets is best***

# Conclusions

Bag of words and deep net profit from each other


Encoding Fisher with FLAIR is 18x faster

Area independent

Supports spatial pyramids, power norm, L2 norm

No approximation

Allows for large-scale localization in video